# Using a Hugging Face Model with Llama.CPP

#### **Docs**

https://github.com/ggerganov/llama.cpp

# **PYTHON WRAPPER (my holy grail for docs now):**

https://github.com/abetlen/llama-cpp-python https://llama-cpp-python.readthedocs.io/en/latest/

# Listen these LLM folk are more MLM, jk. But you do need a GGUF file.

This video kind of helped a lot?

Helped piece the idea that I really need to quantize it or get a gguf file.

<a href="https://www.youtube.com/watch?v=jOEu0PE4ozM">https://www.youtube.com/watch?v=jOEu0PE4ozM</a>

# step 1 <a href="https://github.com/abetlen/llama-cpp-python?">https://github.com/abetlen/llama-cpp-python?</a> <a href="mailto:tab=readme-ov-file#pulling-models-from-hugging-face-hub">tab=readme-ov-file#pulling-models-from-hugging-face-hub</a>

Pull this stuff,

### step 1.1:

https://huggingface.co/docs/huggingface\_hub/en/guides/cli

#### LOGIN AND STUFF

https://huggingface.co/docs/huggingface hub/en/guides/cli#huggingface-cli-login

#### download the model

https://huggingface.co/docs/huggingface hub/en/guides/cli#huggingface-cli-download

I downloaded the whole thing (REFERENCE VIDEO AT THE TOP, WAS HELPFUL ALSO A LOT OF SOURCES TELL YOU IT'S USUALLY FOUND IN THE ~CACHE) <a href="https://huggingface.co/meta-llama/Meta-Llama-3-8B">https://huggingface.co/meta-llama/Meta-Llama-3-8B</a>

#### Once you have this download you can move onto llama.cpp

video helps guide one to getting to this file if you need it

# step 2 <a href="https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#build">https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#build</a>

Here I made a new dir

mkdir IlamaCPP

git clone > clone repo

cd into repo,

https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#cuda

```
$ https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#cuda
```

#### after was able to run stuff here: step3

```
# obtain the official LLaMA model weights and place them in ./models
ls ./models
llama-2-7b tokenizer_checklist.chk tokenizer.model
# [Optional] for models using BPE tokenizers
ls ./models
<folder containing weights and tokenizer json> vocab.json
# [Optional] for PyTorch .bin models like Mistral-7B
ls ./models
<folder containing weights and tokenizer json>
# install Python dependencies
python3 -m pip install -r requirements.txt
# convert the model to ggml FP16 format
python3 convert-hf-to-gguf.py models/mymodel/
# quantize the model to 4-bits (using Q4_K_M method)
./llama-quantize ./models/mymodel/ggml-model-f16.gguf ./models/mymodel/ggml-model-
Q4_K_M.gguf Q4_K_M
# update the gguf filetype to current version if older version is now unsupported
./llama-quantize ./models/mymodel/ggml-model-Q4_K_M.gguf ./models/mymodel/ggml-
model-Q4_K_M-v2.gguf COPY
```

# Step 3 <u>https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#prepare-and-quantize</u>

Quantize the model,

**DYOR** 

# all random stuff i've looked up

#### Not bad resources

https://medium.com/@ingridwickstevens/quantization-of-llms-with-llama-cpp-9bbf59deda35

Like I didn't do this, there's a lot of research, don't overwhelm yourself

https://github.com/ggerganov/llama.cpp/issues/1344

#### my specs,

https://www.google.com/search?

q=can+i+run+llama+3+8b+with+a+3080+site:www.reddit.com&sca\_esv=ba8d56c099ffe98e&sx srf=ADLYWIKY-tjMwps6h-XzWzleX0won6-

<u>28g:1718280528125&sa=X&ved=2ahUKEwjokIPkxdiGAxVDSzABHctnCy4QrQloBHoECCYQBQ&biw=1920&bih=919&dpr=1#ip=1</u>

### read about q6\_k

https://www.google.com/search?

q=quantize+the+model+to+Q6\_K&sca\_esv=ba8d56c099ffe98e&sxsrf=ADLYWILi4Ogav6biasP G7wlhTrkWVYaszw%3A1718280380246&ei=vOBqZr3cDq\_qkvQPx\_qL6Ag&ved=0ahUKEwi9q sGdxdiGAxUvtYQIHUf9Ao0Q4dUDCBA&uact=5&oq=quantize+the+model+to+Q6\_K&gs\_lp=Eg xnd3Mtd2l6LXNlcnAiGnF1YW50aXplIHRoZSBtb2RlbCB0byBRNl9LMgUQIRigATIFECEYoAEy BRAhGKABMgUQIRigATIFECEYoAFItcUBUJYCWJXBAXAAeAKQAQCYAWigAcsDqgEDNC4x uAEDyAEA-AEB-

<u>AECmAlGoALeA8ICBBAAGEfCAgUQABiABMICBhAAGBYYHsICCxAAGIAEGIYDGIoFwgIIEA</u> <u>AYgAQYogTCAggQABgWGAoYHpgDAlgGAZAGCJIHAzUuMaAHqB8&sclient=gws-wiz-serp</u>

### reddits

- 1. <a href="https://www.reddit.com/r/LocalLLaMA/comments/183ie9t/ctransformers\_vs\_llamacpppytho">https://www.reddit.com/r/LocalLLaMA/comments/183ie9t/ctransformers\_vs\_llamacpppytho</a>
  <a href="mailto:n\_which\_one\_should/">n\_which\_one\_should/</a>
- https://www.reddit.com/r/LocalLLaMA/comments/1c8u0n5/thanks\_zuckmeta\_for\_these\_gr eat\_llama\_3\_models\_3/
- 3. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1aeftyu/upgrading\_gtx\_1060\_6gb\_to\_rtx">https://www.reddit.com/r/LocalLLaMA/comments/1aeftyu/upgrading\_gtx\_1060\_6gb\_to\_rtx</a>
  \_3070\_ti\_8gb\_is\_good/

- 4. <a href="https://www.reddit.com/r/LocalLLaMA/comments/148geaj/utilize\_my\_current\_hardware\_or\_upgrade/">https://www.reddit.com/r/LocalLLaMA/comments/148geaj/utilize\_my\_current\_hardware\_or\_upgrade/</a>
- 5. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1c9v4u3/llama\_3\_8b\_any\_way\_to\_run\_t">https://www.reddit.com/r/LocalLLaMA/comments/1c9v4u3/llama\_3\_8b\_any\_way\_to\_run\_t</a> <a href="https://www.reddit.com/r/LocalLLaMA/comments/1c9v4u3/llama\_3\_8b\_any\_way\_to\_run\_t">https://www.reddit.com/r/LocalLLaMA/comments/1c9v4u3/llama\_3\_8b\_any\_way\_to\_run\_t</a>
- 6. USE Q6\_K WHENEVER?

  <a href="https://www.reddit.com/r/LocalLLaMA/comments/1c8mvmc/llama\_3\_discussion\_about\_qu">https://www.reddit.com/r/LocalLLaMA/comments/1c8mvmc/llama\_3\_discussion\_about\_qu</a>
  <a href="mailto:antibarea">ants\_performance/</a>
- 7. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1cci5w6/quantizing\_llama\_3\_8b\_seems\_more\_harmful\_compared/">https://www.reddit.com/r/LocalLLaMA/comments/1cci5w6/quantizing\_llama\_3\_8b\_seems\_more\_harmful\_compared/</a>
- 8. <a href="https://www.reddit.com/r/LocalLLaMA/comments/144uc0l/damn\_i\_was\_so\_satisfied\_with\_my\_3080\_with\_10gb\_of/">https://www.reddit.com/r/LocalLLaMA/comments/144uc0l/damn\_i\_was\_so\_satisfied\_with\_my\_3080\_with\_10gb\_of/</a>
- 9. <a href="https://www.reddit.com/r/SillyTavernAl/comments/1cu4dd6/can\_anyone\_recommend\_som\_e\_local\_models\_for\_3080/">https://www.reddit.com/r/SillyTavernAl/comments/1cu4dd6/can\_anyone\_recommend\_som\_e\_local\_models\_for\_3080/</a>
- 10. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1aezi29/difference\_between\_the\_di
- 11. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1cetn9z/quantization\_seems\_to\_hurt\_the\_quality\_of\_llama\_3/">https://www.reddit.com/r/LocalLLaMA/comments/1cetn9z/quantization\_seems\_to\_hurt\_the\_quality\_of\_llama\_3/</a>
- 12. <a href="https://www.reddit.com/r/LocalLLaMA/comments/1c9qufe/note\_on\_llama\_3\_quantized\_mo">https://www.reddit.com/r/LocalLLaMA/comments/1c9qufe/note\_on\_llama\_3\_quantized\_mo</a> dels/

# https://huggingface.co/docs/huggingface\_hub/en/guides/cli

# ITS RUNNING ON MY CPU

## TO CUDA WE GO

```
pip install llama-cpp-python --extra-index-url https://abetlen.github.io/llama-cpp-
python/whl/cu124 --upgrade --force-reinstall --no-cache-dir
```

#### useful commands

https://stackoverflow.com/questions/78415856/detecting-gpu-availability-in-llama-cpp-python

# watch ur nvidia card

watch -n 1 nvidia-smi

#### takes a screenshot:

```
nvidia-smi
```

#### check cuda version

```
nvcc --version
```

### commands 1 by 1

https://developer.nvidia.com/cuda-12-4-0-download-archive? target\_os=Linux&target\_arch=x86\_64&Distribution=WSL-Ubuntu&target\_version=2.0&target\_type=deb\_local

https://developer.nvidia.com/cuda-12-4-0-download-archive?
target\_os=Linux&target\_arch=x86\_64&Distribution=Ubuntu&target\_version=20.04&target\_type
=deb\_local

https://www.reddit.com/r/KoboldAl/comments/16op2jv/can\_someone\_eli5\_how\_to\_calculate\_the\_number\_of/

https://docs.nvidia.com/cuda/cuda-installation-guide-linux/index.html#wsl https://docs.nvidia.com/cuda/cuda-installation-guide-linux/index.html#post-installation-actions

```
{\tt CMAKE\_ARGS="-DLLAMA\_CUDA=on"} \  \, {\tt pip} \  \, {\tt install} \  \, {\tt llama-cpp-python}
```

```
CMAKE_ARGS="-DLLAMA_BLAS=ON -DLLAMA_BLAS_VENDOR=OpenBLAS" \
    pip install llama-cpp-python
```

#### FINAL COMMAND

```
pip install llama-cpp-python --extra-index-url https://abetlen.github.io/llama-cpp-
python/whl/cu124 --upgrade --force-reinstall --no-cache-dir
```

https://github.com/ggerganov/llama.cpp/issues/6898

checking for cuda existence

```
mako77@DESKTOP-UMFPB6H:~$ ls /usr/local/cuda
DOCS
         README compute-sanitizer extras include libnvvp
                                                                    nvml share
targets version.json
EULA.txt bin doc
                                                   nsightee plugins nvvm src
                                   gds lib64
tools
mako77@DESKTOP-UMFPB6H:~$ ls /usr/local/
bin cuda cuda-12 cuda-12.4 cuda-12.5 etc games include lib man sbin
share src
mako77@DESKTOP-UMFPB6H:~$ export CUDA HOME=/usr/local/cuda-12.4
mako77@DESKTOP-UMFPB6H:~$ which nvcc
/usr/local/cuda-12.4/bin/nvcc
mako77@DESKTOP-UMFPB6H:~$ echo $CUDA HOME
/usr/local/cuda-12.4
```

# **SETTING UP LLAMA CPP**

# **HANDLING THE CPU BEING 100%**

# YOU NEED TO PAY ATTENTION TO ENVINRONMENT VARIABLES, LEARN TO PROPERLY USE THEM, SET THEM, FIND THEM, ETC.

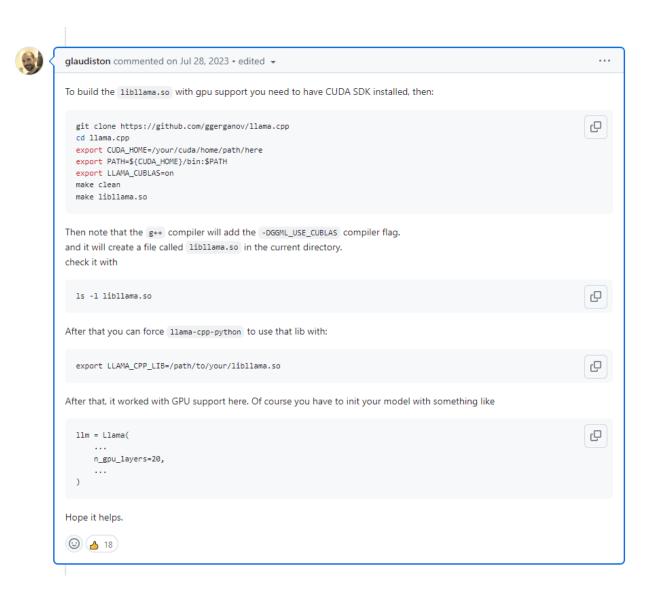
I do believe this is one of the things that saved me

```
CMAKE_ARGS="-DLLAMA_CUBLAS=on" pip install llama-cpp-python --force-reinstall --upgrade --no-cache-dir
```

```
# Linux and Mac
CMAKE_ARGS="-DLLAMA_BLAS=ON -DLLAMA_BLAS_VENDOR=OpenBLAS" \
   pip install llama-cpp-python
```

#### **Resources for CPU BEING 100%**

(1) Rebuilding, at first I did not build with cuda home or paths and stuff set. https://github.com/abetlen/llama-cpp-python/issues/509#issuecomment-1655864538



After doing all that,

ALSO MAKING SURE I INSTALLED CUDA, BECAUSE I DID NOT KNOW I WOULD NEED CUDA LOL

I re-quantized the model after this, meaning I deleted the previous one and did it with these conditions fulfilled.

## After re-quantizing

```
# Linux and Mac
CMAKE_ARGS="-DLLAMA_BLAS=ON -DLLAMA_BLAS_VENDOR=OpenBLAS" \
   pip install llama-cpp-python
```

I guess it did say to do this prior,

The official docs of this wrapper suggest this:

# **Upgrading and Reinstalling**

To upgrade and rebuild llama-cpp-python add --upgrade --force-reinstall --no-cachedir flags to the pip install command to ensure the package is rebuilt from source.

Which was almost right,

but ultimately what worked for me was plainly making sure I pip uninstalled the package

then ran

```
CMAKE_ARGS="-DLLAMA_CUBLAS=on" pip install llama-cpp-python --force-reinstall --upgrade --no-cache-dir
```

Links that helped ---

#### found solution here

https://github.com/abetlen/llama-cpp-python/issues/576#issuecomment-1766003289 https://github.com/abetlen/llama-cpp-python/issues/576

huge helped me with the idea to rebuild everything. When in doubt just rebuild

https://stackoverflow.com/questions/78415856/detecting-gpu-availability-in-llama-cpp-python Definitely need to set env variables first i imagine

# I did do this also when running the build

https://github.com/ggerganov/llama.cpp?tab=readme-ov-file#cuda really feel like it didnt do much, but who knows i did it so its included

# this is where i found the command

https://github.com/ggerganov/llama.cpp/issues/6360